

CAFA: a Controllable Automatic Foley Artist

Anonymous ICCV submission

Paper ID 5746

Abstract

001 Foley is a key element in video production, refers to the process of adding an audio signal to a silent video while ensuring semantic and temporal alignment. In recent years, the rise of personalized content creation and advancements in automatic video-to-audio models have increased the demand for greater user control in the process. One possible approach is to incorporate text to guide audio generation. While supported by existing methods, challenges remain in ensuring compatibility between modalities, particularly when the text introduces additional information or contradicts the sounds naturally inferred from the visuals. In this work, we introduce CAFA (Controllable Automatic Foley Artist) a video-and-text-to-audio model that generates semantically and temporally aligned audio for a given video, guided by text input. CAFA is built upon a text-to-audio model and integrates video information through a modality adapter mechanism. By incorporating text, users can refine semantic details and introduce creative variations, guiding the audio synthesis beyond the expected video contextual cues. Experiments show that besides its superior quality in terms of semantic alignment and audio-visual synchronization the proposed method enable high textual controllability as demonstrated in subjective and objective evaluations. ¹

025 1. Introduction

026 In recent years, personal content creation has become a major part of everyday life, shaping how we work, entertain, 027 and communicate. One example is *Foley*, the art of adding 028 sound effects to silent videos while ensuring precise semantic 029 and temporal alignment [1]. Traditionally, this process 030 was done manually by professional sound designers. How- 031 ever, with the growing demand for fast and immediate per- 032 sonal digital content, the need for automation and access- 033 ability of this process has increased. An effective Foley 034 generation approach should produce high-quality, synchro-

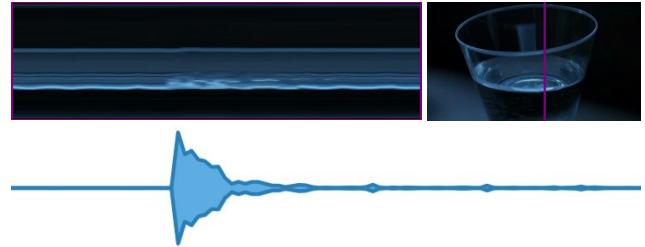


Figure 1. **Motivation.** An iconic scene from *Jurassic Park*, where water in a glass shakes due to the approaching footsteps of a T-Rex. Inferring the generated sound from the video alone is insufficient, as the task is inherently ambiguous. **Top:** a representative frame and a Y-T slice (from the purple column), where the temporal cue of the shake is faintly visible. **Bottom:** Our method leverages the prompt "T-Rex Stomping" to generate a synchronized audio track that aligns with both the visual timing and artistic intent.

nized audio while also allowing users to creatively shape the sound, balancing precision with creative flexibility.

Building on this need, recent advancements in generative models have led to the development of Video-to-Audio (V2A) models, which aim to automate Foley synthesis and explore cross-modal correspondences [17, 31, 43, 46, 48, 50]. While these models effectively capture semantic information through global visual representations such as CLIP [36], they often rely on motion-sound relationships for temporal alignment [31, 46, 48]. Some approaches, including [17], model motion explicitly using optical flow [14], while others [31, 46, 48] leverage contrastive learning-based encoders such as CAVP [31] and AV-CLIP [19] to learn temporally and semantically aligned audio-visual features. Despite these advancements, existing models remain limited to extracting information from the video itself and struggle to incorporate user-provided cues, restricting flexibility and creative control over sound design.

To bridge this gap, Text-and-Video-To-Audio (TV2A) models have been introduced, integrating textual information to enhance control over audio generation [5, 7, 20, 28, 50, 55]. By incorporating text, these models allow users to modify audio semantics, add details, and generate diverse sound variations. For instance, text can specify how

¹Samples and code can be found in our anonymized [demo page](#)

060 a sound should be perceived, such as describing a door as
 061 creaking or coffee being sipped loudly. Another possibility
 062 is introducing creativity through text; a barking dog in
 063 a video could instead sound like a meowing cat or a crow-
 064 ing rooster, depending on the accompanying description. In
 065 the context of soundtrack design, one would often like to
 066 add sounds which do not appear in the video, such as in
 067 the iconic scene from *Jurassic Park*, where water in a glass
 068 shakes due to the approaching footsteps of a T-Rex; see Figure
 069 [1] for a visual example. However, textual conditioning
 070 is often not sufficiently strong or may come at the expense
 071 of temporal alignment between video and audio. Additionally,
 072 when the text describes semantics that differ from the
 073 video, existing models frequently struggle to generate a nat-
 074 ural and coherent audio signal (See Section 5).

075 Various methods have been explored for integrating text
 076 into multimodal systems. A common strategy involves
 077 jointly training video, text, and audio representations to
 078 capture shared semantics [5, 7]. However, this requires
 079 retraining the entire network whenever modifications are
 080 made, leading to high computational costs. Alternatively,
 081 a training-free method [50] leverages a shared latent space
 082 to link the modalities, eliminating the need for retraining.
 083 Yet, this introduces test-time optimization, increasing in-
 084 ference time and potentially degrading output quality and
 085 alignment. A middle-ground solution employs a modali-
 086 ty adapter (e.g., ControlNet mechanism [54]), which uses
 087 video inputs to condition a pretrained Text-to-Audio (T2A)
 088 model [20, 55], providing an effective way to incorporate
 089 video information into text-driven audio synthesis.

090 In this work, we introduce CAFA, which stands for Con-
 091 trollable Automatic Foley Artist, a novel text-and-video-to-
 092 audio model that extends beyond temporal and semantic
 093 synchronization, allowing users to shape and control sound
 094 through textual cues. CAFA leverages a ControlNet like
 095 modality adapter to flexibly integrate pretrained T2A mod-
 096 els with video-based features while maintaining a relatively
 097 low training cost (48 A100 GPU hours for CAFA vs. 304
 098 H100 GPU hours for the baseline method). Specifically, we
 099 explore Stable-Audio-Open [9] and TangoFlux [16] as T2A
 100 models. To extract temporal and semantic features, we ex-
 101 periment with AVCLIP [19] and CLIP [36] as the video rep-
 102 resentations. CAFA achieves high-quality audio synthesis,
 103 temporal synchronization, and contextual alignment perfor-
 104 mance comparable to state-of-the-art V2A and TV2A mod-
 105 els. Additionally, it significantly surpasses existing TV2A
 106 approaches when the text and visual conditioning cues are
 107 semantically different, demonstrating greater control over
 108 generated sound.

109 Our main contributions are: (i) We introduce CAFA, a
 110 novel TV2A model that allows the generation of tempo-
 111 rally and semantically aligned audio while providing ex-
 112 tensive textual control over the generated audio; (ii) We

113 evaluate CAFA against existing V2A and TV2A mod-
 114 els, demonstrating comparable performance in audio qual-
 115 ity and video-audio compatibility, while achieving superior
 116 performance for textual control, as validated through disen-
 117 tanglement experiments, objective evaluations, and human
 118 studies; (iii) CAFA is built on the modality adaptation (via a
 119 ControlNet mechanism), enabling precise temporal control
 120 while offering a versatile framework that supports modular
 121 integration, accommodating different T2A models (Stable
 122 Audio Open and TangoFlux). Additionally, it facilitates the
 123 efficient incorporation of video representations, leading to
 124 more effective training compared to alternative methods.

2. Background

2.1. Latent Diffusion Models

125 Latent Diffusion Models (LDMs) [39] are a class of gener-
 126 ative models that perform a diffusion process within a
 127 learned latent space \mathbf{z} . A Variational Autoencoder (VAE)
 128 encodes a data sample $\mathbf{x} \sim p(\mathbf{x})$ into a lower-dimensional
 129 latent space $\mathbf{z} \in \mathbb{R}^d$ using an encoder \mathcal{E} , while a decoder \mathcal{D}
 130 reconstructs \mathbf{x} . Performing diffusion in this reduced space
 131 significantly reduces computational cost while maintaining
 132 high-quality generation.

133 The diffusion process follows two Markovian paths: the
 134 forward and reverse processes. In the forward process, a
 135 clean latent representation \mathbf{z}_0 is gradually corrupted with
 136 additive Gaussian noise:

$$\mathbf{z}_t = \sqrt{\alpha_t} \mathbf{z}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_t, \quad (1)$$

139 where $\{\alpha_t\}_{t=1}^T$ defines the noise schedule, $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 140 and $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. A key consequence of the forward pro-
 141 cess is its marginal distribution:

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2)$$

144 where $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$. A neural network, trained as a de-
 145 noiser, learns to estimate ϵ given the noisy input \mathbf{z}_t , the
 146 timestep t , and conditioning information c , such as en-
 147 coded text. The training objective minimizes the difference
 148 between the true noise and the predicted noise:

$$\mathcal{L} = \mathbb{E}_{\mathbf{z}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t, c} [\|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, c)\|]. \quad (3)$$

149 By using the network output, the reverse process aims to
 150 reconstructs \mathbf{z}_0 from \mathbf{z}_T by iteratively denoising it. While
 151 initial works [13, 44] formulated this process discretely,
 152 Song et al. [45] showed that it can be equivalently expressed
 153 as an Ordinary Differential Equation (ODE) which can be
 154 solved numerically using dedicated solvers. Specifically,
 155 Stable Audio Open [9] employs DPM-Solver++ [30] and
 156 follows the the v-objective approach [41].

157 **Classifier-Free-Guidance (CFG)** is a widely used method
 158 to improve performance in conditional generative models,

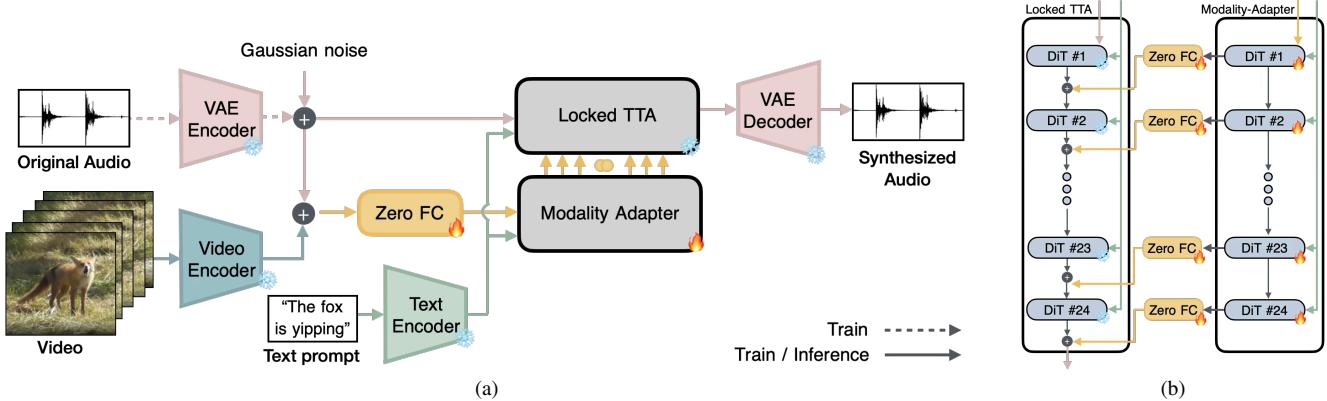


Figure 2. (a) **Method overview:** our model is text-and-video-to-audio, leverages pretrained models for audio generation, and video encoding. The original audio and VAE audio encoder are only used during training. (b) **Adaptor:** Illustration of the internal connectivity between the backbone T2A model and our video conditioning adaptor, with fully connected (FC) layers explicitly shown.

originally demonstrated in diffusion-based image generation approaches [12]. CFG is an effective control mechanism for steering the inference process to better align with provided conditioning signals. Specifically, it modifies the predicted noise by linearly combining the estimates from a conditional diffusion model and a jointly trained unconditional model, resulting in the following formulation:

$$\tilde{\epsilon}_\theta(\mathbf{z}_t, c, t) = \epsilon_\theta(\mathbf{z}_t, t) + \gamma (\epsilon_\theta(\mathbf{z}_t, t, c) - \epsilon_\theta(\mathbf{z}_t, t)) \quad (4)$$

where γ determines the strength of the guidance, with higher values enforcing stronger adherence to the conditioning signal.

2.2. ControlNet Mechanism

The ControlNet mechanism was initially introduced as a neural network architecture for controlling text-to-image models through spatially localized conditioning (e.g., canny edge and depth maps) [54]. It preserves the quality and stability of a large pretrained model by locking its weights, while enabling the incorporation of control signals through a replicated copy of that backbone model. These components are connected via zero-initialized convolutional layers, allowing a gradual integration which minimizes noise interference during training.

3. Method

CAFA consists of two main components: a pretrained backbone Text-To-Audio (T2A) model, with frozen weights to maintain audio quality, and a trainable modality adapter that integrates temporal and semantic video information. The components are linked through Zero Fully-Connected (FC) Layers, which prevent noise from disrupting the backbone model during early training. This structure allows us to benefit from the long pre-training of the foundational

T2A model instead of training all three modalities from scratch. Figure 2a provides a high-level overview of the proposed method.

Text-to-Audio Backbone. A variational autoencoder [9] encodes the input signal $\mathbf{x} \in \mathbf{R}^{2 \times L}$ (2 channels for stereo) into a latent representation $\mathbf{z} \in \mathbf{R}^{T \times C}$, with L denoting the temporal length of the audio, while T and the $C = 64$ correspond to the temporal dimension and feature size, respectively. Next, noise is added to the latent representation, producing \mathbf{z}_t , which is then processed by a core architecture built from a stack of Diffusion Transformer (DiT) blocks [35], with model-dependent variations. Each DiT block is controlled by a text input, encoded by a pretrained text encoder, guiding the generation process. Furthermore, a timing mechanism sets the signal length and fills the rest with silence.

We experiment with two T2A models: Stable Audio Open [9] and TangoFlux [16]. These models provide the underlying architecture for creating high-quality stereo audio content at a sampling rate of 44.1 kHz based on textual descriptions. Although the architectural designs and sampling methods differ, both models follow a comparable core structure. This allows us to evaluate the flexibility of our approach, demonstrating that it is not tied to a single family of models and can be applied across various T2A models.

Modality Adapter. The adapter is tasked with incorporating the video to guide the T2A model generation. Our design is inspired by ControlNet [54], with some notable differences. First, it operates on the temporal domain, requires synchronization of features from different modalities. Second, our T2A model is a DiT, rather than a U-Net [40], changing the connectivity between the adapter and the base model. Specifically, after the preprocessing stage, informative features from the video, $E_v \in \mathcal{R}^{T \times C}$, are passed through a Zero FC layer and added to \mathbf{z}_t . Additionally, the

226 hidden states, extracted from each DiT block, are processed
 227 through Zero FC layers and added to the backbone model,
 228 as depicted in Figure 2b. Only the adapter is being trained,
 229 while the T2A, Text Encoder, and Video Encoder, all kept
 230 frozen.

231 **Video Representation.** We experiment with AVCLIP [19]
 232 and CLIP [36], both of which are trained with contrastive
 233 learning. AVCLIP processes 0.64-second video-audio seg-
 234 ments and applies InfoNCE loss [33] to differentiate be-
 235 tween positive and negative samples. Its video encoder,
 236 built on MotionFormer [34], enhances the modeling of dy-
 237 namic scenes by capturing implicit motion paths. To ensure
 238 better alignment with \mathbf{z}_t , we increased the segment overlap
 239 compared to the original work, where 216 samples roughly
 240 correspond to 10 seconds. We further examine the contribu-
 241 tion of CLIP to the semantic scene understanding, similarly
 242 to MMAudio [7]. Since CLIP is less sensitive to object lo-
 243 cations [3], we introduce a temporal dimension to capture
 244 significant visual changes. We extract CLIP representations
 245 at 5 FPS and interpolate to match the temporal dimension
 246 of \mathbf{z}_t . Finally, both representations are preprocessed before
 247 being integrated with the modality adapter, as detailed in
 248 the supplementary [25].

249 **Asymmetric Classifier-Free Guidance.** Incorporating
 250 CFG into our model with text as the conditioning signal c
 251 presents challenges in effectively utilizing video condition-
 252 ing at higher guidance scales. To address this limitation,
 253 we propose *Asymmetric Classifier-Free Guidance*, where
 254 the modality adapter’s output is selectively modulated in
 255 the conditional and unconditional pathways during syn-
 256 thesis. Unlike the standard approach, which equally integrates
 257 the modality adapter into the backbone model in both path-
 258 ways, our method introduces an asymmetric scaling factor,
 259 $0 \leq \alpha \leq 1$, reducing the influence of the adapter in the
 260 unconditional path,

$$h'_{i,c} = h_{i,c}, \quad h'_{i,uc} = \alpha \cdot h_{i,uc}, \quad (5)$$

261 where $h_i = [h_{i,c}, h_{i,uc}]$ are the hidden states, with $h_{i,c}$ rep-
 262 resenting the conditional path and $h_{i,uc}$ the unconditional
 263 path. Consequently, under $\alpha < 1$, this adjustment induced
 264 controlled disparity between $\epsilon_\theta(\mathbf{z}_t, t, c)$ and $\epsilon_\theta(\mathbf{z}_t, t)$ effec-
 265 tively amplifies the video conditioning signal. Standard
 266 CFG is a special case, corresponding to $\alpha = 1$. Our exper-
 267 iments demonstrate that this simple yet effective modifica-
 268 tion significantly enhances adherence to video conditioning
 269 while maintaining high generation quality and text control-
 270 lability.

272 4. Experimental Setup

273 4.1. Datasets

274 The proposed model was trained using two datasets: VG-
 275 GSound [2] and VisualSound [46]. VisualSound is a sub-

276 set of VGGSound filtered to include samples with high
 277 ImageBind [11] scores. Both datasets contain 10-second
 278 video clips across diverse acoustic categories accompanied
 279 by video captions. For TV2A and V2A evaluation, we use
 280 the VGGSound dataset and VGGSound-Sparse [18], which
 281 is a subset of VGGSound containing 12 categories of natu-
 282 rally sparse audio events such as “dog barking” or “playing
 283 tennis”.

284 4.2. Baseline Methods

285 We compare CAFA against several state-of-the-art mod-
 286 els, namely MMAudio [7] (*large_44k_v2* version), Foley-
 287 Crafter [55], VATT [28], ReWaS [20], Frieren [48], and
 288 MultiFoley [6]. For FoleyCrafter, we follow original con-
 289 figurations by formatting text prompts as “The sound of
 290 <label>”. For VATT, Frieren, and MultiFoley, we con-
 291 sider the samples provided by the respective authors. Mul-
 292 tiFoley samples were only available for a subset of the test
 293 set with high ImageBind scores. ReWaS is evaluated us-
 294 ing its default configuration. Notice, with 5-second sam-
 295 ples, unlike other models that produce 8-second samples.
 296 Hence, for a fair comparison against this model we truncate
 297 the videos to 5-second.

298 4.3. Implementation Details

299 CAFA models are initially trained for 48k steps on VG-
 300 GSound, followed by fine-tuning for 33k more steps on Vi-
 301 sualSound. Training was performed with a batch size of 16,
 302 using the AdamW [29] optimizer, on a single A100 GPU.
 303 We generate samples using $CFG = 7$, Asymmetric CFG
 304 scale of $\alpha = 0.5$, and 50 inference steps, while keeping the
 305 rest of the TTA model configuration unchanged. Our model
 306 is trained on 10-second samples, and the output truncated to
 307 8 seconds for fair comparison with the baseline methods.

308 4.4. Evaluation Metrics

309 We evaluate model performance across four complementary
 310 dimensions that capture different aspects of audio-visual
 311 generation: Audio Quality, Audio-Visual Semantic Align-
 312 ment, Audio-Visual Temporal Alignment, and Audio-Text
 313 Semantic Alignment.

314 **Audio Quality.** We employ three established metrics to
 315 assess the fidelity and naturalness of generated audio: (i)
 316 Fréchet Audio Distance (FAD) [21] which measures distri-
 317 butional similarity between features extracted from ground
 318 truth and generated audio; (ii) Kullback-Leibler Distance
 319 (KL) [24], which quantifies the difference between prob-
 320 ability distributions of per-sample ground truth and gen-
 321 erated audio features; and (iii) Inception Score (IS) [42]
 322 that evaluates the generated audio quality independently of
 323 ground truth references. We utilize PANNS [22] as the fea-
 324 tures extractor for all three audio quality metrics.

Model	FAD↓	IS↑	CLAP↑	Acc↑	DeSync↓
FC	57.00	<u>6.10</u>	0.10	0.69	1.30
MMA	16.43	6.77	0.10	0.36	0.57
ReWaS	38.94	5.13	0.09	0.74	1.19
CAFA (Ours)	<u>27.33</u>	5.63	0.21	0.87	0.81

Table 1. **Semantically different text and video conditioning.** Our method surpasses strong concurrent SOTA in terms of prompt adherence by a large margin. Arrows indicate whether higher (\uparrow) or lower (\downarrow) values are better. FC: FoleyCrafter, MMA: MMAudio.

325 **Audio-Visual Semantic Alignment.** We leverage Image-
 326 Bind (IB) [11] to quantify semantic similarity between the
 327 ground truth video and generated audio. This cross-modal
 328 embedding model measures whether the generated audio
 329 contains appropriate sounds for the visual content.

330 **Audio-Visual Temporal Alignment.** We utilize DeSync [7] (also known as Sync [46] or AV-Sync [6])
 331 to measure temporal synchronization between audio and
 332 video. DeSync calculates the average absolute offset (in
 333 seconds) between ground truth video and generated audio
 334 using Synchformer [19] predictions. Following prior
 335 work [7], we average DeSync scores from the first and
 336 last 4.8 seconds of audio to accommodate Synchformer’s
 337 limited context window.

339 **Audio-Text Semantic Alignment.** We employ CLAP [49]
 340 to evaluate similarity between generated audio and textual
 341 descriptions of the video by calculating the cosine similarity
 342 between the text and audio embeddings.

343 5. Results

344 We start by evaluating model performance when considering
 345 semantically different text and video conditioning. Ideally,
 346 we expect the model to generate textually described au-
 347 dio aligned with the visual cues. For that we generate audio
 348 for each video in VGGSound-Sparse, using the ground truth
 349 video paired with captions from each of the 11 other cate-
 350 gories. This cross-category approach creates a challenging
 351 scenario where models must follow textual instructions that
 352 deliberately conflict with visual content.

353 Notice, under this setup we compute the CLAP similar-
 354 ity score between the new caption and generated audio. We
 355 also use CLAP as a classifier between the new caption and
 356 the GT caption for the generated audio, reporting binary
 357 classification accuracy (Acc).

358 For FoleyCrafter, we follow [6] and disable the semantic
 359 adapter to allow the model to generate the requested cap-
 360 tion, using only the temporal adapter to retrieve information
 361 from the video.

362 Results are summarized in Table 1 with visual exam-
 363 ples depicted in Figure 3. MMAudio achieves the high-

Model	FAD↓	KL↓	IS↑	IB↑	DeSync↓	CLAP↑
FC†	13.68	2.56	10.68	<u>0.27</u>	1.30	<u>0.12</u>
MMA†	5.32	1.64	17.18	0.33	0.77	0.23
Frieren†	11.76	2.70	12.33	0.23	1.04	0.11
FC	22.17	2.87	13.30	0.16	1.31	0.18
MMA	6.89	<u>1.65</u>	20.44	0.34	0.76	0.25
VATT	<u>11.13</u>	1.48	11.85	0.25	1.28	0.15
ReWaS*	14.71	2.69	8.45	0.15	1.18	0.18
MF*	13.51	<u>1.65</u>	<u>15.89</u>	<u>0.27</u>	1.04	<u>0.23</u>
CAFA (Ours)	12.60	2.02	13.45	0.21	<u>0.96</u>	<u>0.23</u>

Table 2. **Quantitative comparison.** We report results comparing standard V2A models, V2A variants of TV2A models (indicated by †), TV2A models, and our method. FC:FoleyCrafter, MMA:MMAudio, MF:MultiFoley. * indicates variations - we compare with ReWaS on samples trimmed to 5 seconds, and compare with MultiFoley on their selected subset of the test set.

364 est audio quality (best FAD and IS scores) and temporal
 365 alignment (DeSync), but critically fails at semantic control.
 366 Its Acc scores of 0.58 (w\negp) and 0.36 (w\o negp)
 367 demonstrate that it systematically generates audio corre-
 368 sponding to the visual content rather than adhering to the
 369 requested text prompt, essentially negating the purpose of
 370 text-guided generation. Despite disabling FoleyCrafter’s
 371 semantic adapter specifically for this experiment, it ex-
 372 hibits severe temporal misalignment (worst DeSync score
 373 of 1.30) producing outputs that fail to synchronize with vi-
 374 sual events. This renders its generation ineffective even
 375 for basic audio-visual correspondence. ReWaS performs
 376 sub-optimally across all evaluation dimensions—producing
 377 lower audio quality, weaker temporal alignment, and poorer
 378 semantic control compared to our approach—without ex-
 379 hibiting a particular strength in any area to compensate for
 380 these deficiencies. In contrast, CAFA successfully balances
 381 all critical requirements: achieving strong audio quality
 382 (second-best FAD and comparable IS), effective temporal
 383 alignment (second-best DeSync), while significantly out-
 384 performing all competitors in semantic controllability with
 385 a CLAP score of 0.21 (compared to < 0.1 for all others)
 386 and Acc of 0.87. This comprehensive performance profile
 387 makes our approach uniquely capable of producing high-
 388 fidelity, temporally-aligned audio that accurately follows
 389 semantic text instructions.

390 **Semantically Aligned Visual and Textual Conditions.**
 391 Next, we compare CAFA where the visual and textual con-
 392 ditions are semantically aligned. We compare the proposed
 393 method against V2A and VT2A methods, considering the
 394 VGGSound test set using the standard configurations. For
 395 V2A models, we do not use textual descriptions. Results
 396 are presented in Table 2.

397 While MMAudio emerges as the strongest performer,
 398 comparison of MMAudio with and without text conditions
 399 reveals minimal benefits from textual inputs. Specifically,

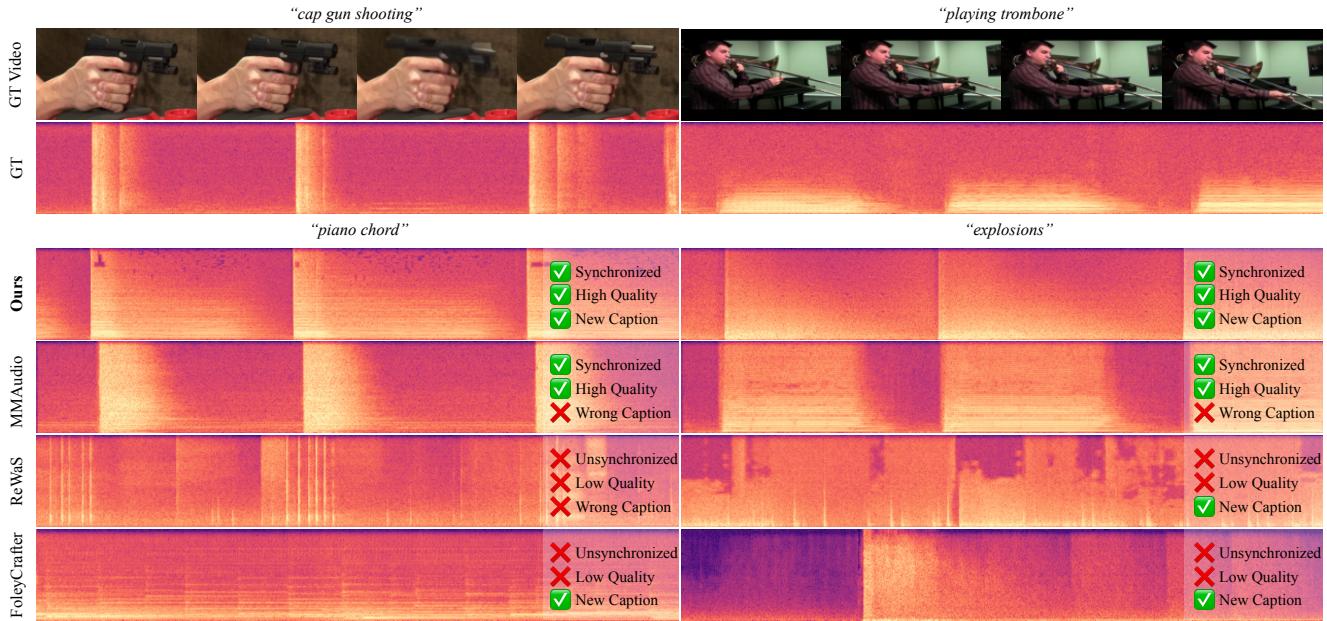


Figure 3. **Qualitative Comparison of Text-Video Disentanglement.** A comparative analysis of various TV2A models: Ground Truth (GT), CAFA (ours), MMAudio, ReWaS, and FoleyCrafter, using the same configurations as in Table I. Our model consistently delivers synchronized, high-quality generations that accurately adhere to the requested target captions, outperforming other approaches. Full videos presented at [our demo page](#).

FAD got slightly worse when text is added, while other metrics remain largely unchanged. Interestingly, a similar pattern emerges in FoleyCrafter, which shows mixed performance with text conditioning - its FAD, KL, IB got worsen, while IS and CLAP improved, and DeSync stays unchanged. These results strengthens our findings in Table I, demonstrating the improved flexibility and controllability of the proposed method compared to the baseline methods.

CAFA demonstrates balanced performance across all evaluation metrics, achieving the second-best scores in both DeSync and CLAP while maintaining competitive audio quality metrics. These results indicate effective multi-modal conditioning without compromising performance in any single dimension.

Human Study. Lastly, we evaluate the subjective quality of audio generated by the proposed method. We conduct a user study using the following evaluation protocol considering both the original textual captions (OC) and visually unaligned text caption (UC). We follow the above mention protocol.

Each participant was presented with a pair of videos side by side. For every pair, the participant was asked three distinct questions: (i) Prompt Adherence (PA): *Which audio better matches the description '[text prompt]'?* This question measured how well the audio corresponds to the prompt; (ii) Alignment: *Which audio better aligns with the*

Comparison	Criterion	% (OC)	% (UC)
CAFA vs MMAudio	Align.	0.24	0.24
	Quali.	0.29	0.26
	PA	0.39	0.87
CAFA vs FoleyCrafter	Align.	0.83	0.95
	Quali.	0.81	0.88
	PA	0.82	0.69
CAFA vs ReWaS	Align.	0.96	0.96
	Quali.	0.95	0.89
	PA	0.94	0.77

Table 3. **Human study.** Win rate results (%) of CAFA against three baseline methods. Results are reported for: time alignment (Align.), audio quality (Quali.), and prompt adherence (PA), considering the original caption (OC) and unaligned caption (UC).

timing of visual movements and events in the video? This question assessed the synchrony between audio events and the corresponding visual actions; (iii) Quality: *Which audio has higher overall technical quality (considering naturalness, clarity, and lack of artifacts)?* This evaluates the technical fidelity of the audio signal. Participants provided independent responses for each of these questions. In both protocols, the ordering of the video pairs was randomized to mitigate any potential bias. The protocol was implemented using a custom web interface built on Amazon SageMaker Ground Truth. We use 17 videos from VGGSound test set

427
428
429
430
431
432
433
434
435
436
437

Model	Audio Quality			A-V Sem.	A-V Tem.	A-T
	FAD↓	KL↓	IS↑	IB↑	DeSync↓	CLAP↑
CAFA-B	12.57	2.04	11.84	<u>0.21</u>	<u>1.00</u>	0.23
CAFA	12.60	2.02	13.45	0.21	0.96	0.23
CAFA-C	14.44	1.98	14.18	0.22	1.02	0.23
CAFA-TF	19.94	2.16	16.94	0.20	1.12	0.23

Table 4. **Model ablation.** We compare different variants of our model, and argue that our default model CAFA strikes an overall good balance between audio quality, and semantic and temporal alignment to video and prompt. CAFA-B - our model before finetuning on VisualSound. CAFA-C - additionally leverage CLIP visual features. CAFA-TF - use TangoFlux [16] as base T2A model.

438 to generate two audio samples: with an original caption and
439 a new one by every method. Then each pair of CAFA vs
440 competitor (in random order) was shown to 6 raters and the
441 results were averaged across 6 annotations per pair.

442 Win rate results of the proposed method against the eval-
443 uated baselines are presented in Table 3. When comparing
444 to FoleyCrafter and ReWas the proposed method achieves
445 superior performance across all setups. As expected, when
446 comparing to MMAudio, the proposed method reach infe-
447 rior performance considering both audio quality and time
448 alignment, however, under the unaligned text and visual
449 conditioning, CAFA reach significantly better performance
450 than MMAudio.

451 6. Analysis

452 **Model Variations.** We first evaluate several architectural
453 variants of CAFA. Specifically, we consider: (i) CAFA-
454 B, which is identical to CAFA but does not include ad-
455 dditional finetuning phase on VisualSound, allowing us to
456 evaluate training efficiency; (ii) CAFA-C, which leverages
457 both AV-CLIP and CLIP as visual encoders, combined via
458 MLP. This model was trained for 84k steps on VGGSound
459 train split with a batch size of 8 and the same optimizer
460 as CAFA; and (iii) CAFA-TF, that leverages TangoFlux as
461 base T2A model with AV-CLIP visual encode. CAFA-TF
462 was trained for 32k steps on VGGSound train split with
463 batch size 64 and the same optimizer as CAFA. For vari-
464 ants using StableAudio-Open, we used a CFG value of 7.0,
465 asymmetric CFG $\alpha = 0.5$, and 50 inference steps. For
466 CAFA-TF, we used a CFG value of 4.5, asymmetric CFG
467 $\alpha = 0.8$, and 50 steps.

468 Results, presented in Table 4, shows that all variants
469 achieve similar performance across metrics, with minimal
470 differences in FAD, IS, KL, and DeSync. The CAFA-
471 C results demonstrate that adding CLIP visual condition-
472 ing alongside AV-CLIP provides no meaningful benefits,
473 highlighting that the combination of AV-CLIP with text
474 conditioning is sufficient for effective audio generation.

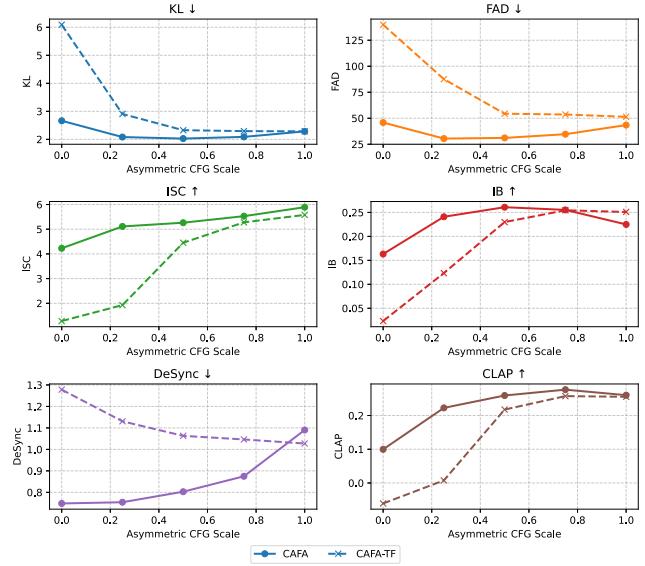


Figure 4. Comparison of Asymmetric CFG Scaling Values. CAFA-TF is our adapter applied with TangoFlux [16], while the default implementation uses StableAudio-Open [9].

The CAFA-B variant, trained for only 48k steps, performs slightly worse than CAFA. Finally, the performance of CAFA-TF confirms that our method generalizes effectively across different T2A base models.

Training Efficiency Analysis. As shown in Table 5 in the Appendix, CAFA demonstrates significant training efficiency compared to other state-of-the-art models. While direct comparison is challenging due to differences in reporting methodologies, hardware configurations, and training approaches, several observations can be made. CAFA requires substantially fewer training steps (81k total) compared to models like Frieren (2.4M steps) and MultiFoley (650k steps). Even our base model (CAFA-B), trained for only 48k steps, achieves performance comparable to models trained for much longer periods. The adapter-based approach allows CAFA to leverage pre-trained text-to-audio models effectively, reducing the need for extensive training from scratch. While MMAudio reports 304 GPU hours on H100 hardware, CAFA estimated 48 GPU hours on A100 hardware represents a more efficient use of compute resources when considering the relative performance. The modular architecture of CAFA facilitates efficient training while maintaining high performance across audio quality, temporal alignment, and textual control metrics. This efficiency analysis further highlights the practical advantages of our approach, making CAFA more accessible for research and potential real-world applications.

Asymmetric CFG Scaling. Finally, we investigate different settings for asymmetric CFG scaling. We evaluate both

504 CAFA and CAFA-TF over the VGGSound-Sparse test set
505 with the same configurations as described at Section 5. Figure
506 4 depicts the results, illustrating the trade-off between
507 different scaling parameters. Our analysis shows that for
508 both models, the trends are similar, with values between 0.5
509 and 1 corresponding to the best results.

510 7. Related work

511 **Text To Audio Models.** The field of Text-to-Audio (T2A)
512 has advanced significantly, with ongoing improvements
513 in text representations and high-quality audio generation.
514 Early approaches included AudioGen [23], an autoregres-
515 sive model employing a discrete waveform representation,
516 while DiffSound [52] adopted discrete diffusion, remov-
517 ing the need for autoregressive token decoding. As the
518 field evolved, models like AudioLDM [26], StableAudio
519 1 [8], and Make-An-Audio [15] leveraged latent diffusion
520 and incorporated CLAP [49] embeddings to improve text
521 decoding. Recognizing the importance of capturing tempo-
522 ral, acoustic, and semantic information, newer models such
523 as Make-An-Audio 2 [15], AudioLDM 2 [27], and Tango
524 [10] integrated large language models (LLMs) [38] to en-
525 hance text-audio alignment. Building on these advan-
526 tages, Tango 2 [32] further refined temporal alignment by
527 employing Direct Preference Optimization (DPO) [37].

528 Our work is built on StableAudio Open [9] and
529 TangoFlux [16], designed for high-quality text-to-audio gen-
530 eration. StableAudio Open is a latent diffusion model that
531 generates stereo audio up to 47 seconds from text input, uti-
532 lizing a T5 text encoder for text processing and enabling
533 control over output length. In contrast, TangoFlux is based
534 on rectified flow, producing stereo audio up to 30 seconds
535 at 44.1 kHz, while leveraging a pretrained autoencoder from
536 StableAudio Open to enhance efficiency. Additionally, Tan-
537 goFlux incorporates CLAP-Ranked Preference Optimiza-
538 tion (CRPO) to generate and refine audio preference data.

539 **Video To Audio Models.** A key step in automating the
540 Foley process is achieved through Video-To-Audio (V2A)
541 models. Early approaches, such as SpecVQGAN [17],
542 RegNet [4], and FoleyGAN [51], used adversarial training
543 and GAN-based architectures to generate high-quality au-
544 dio. Diff-Foley [31], a diffusion-based model, introduced
545 CAVP contrastive learning to improve temporal and seman-
546 tic alignment. Alternatively, Frieren [48], based on recti-
547 fied flow, enables efficient audio generation in fewer steps.
548 V-AURA [46] adopts an autoregressive approach, leverag-
549 ing the AVCLIP [19] representation to extract high-frame-
550 rate temporal and semantic features while bypassing spec-
551 trogram conversion.

552 Beyond models that require training from scratch, V2A-
553 Mapper [47] and Seeing and hearing [50] employ training-
554 free optimization, utilizing pretrained text-to-audio gener-

555 ators or modality mappers to condition audio generation.
556 While these methods reduce computational cost, they often
557 struggle with fine-grained temporal synchronization, high-
558 lighting the ongoing challenge of bridging the gap between
559 video and audio in a seamless and efficient manner.

560 **Text and Video To Audio Models.** Text-and-Video-to-
561 Audio (TV2A) models introduce text conditioning to en-
562 hance control over synthesized audio. VATT [28] leveraged
563 an LLM decoder, functioning as both a video-to-caption
564 model and a video-text-to-audio model. MMAudio [7] and
565 MultiFoley [5] explicitly trained all three modalities from
566 scratch, achieving state-of-the-art results in signal quality
567 and synchronization. While MMAudio introduced a novel
568 network structure for modality fusion, MultiFoley, based on
569 DiT [35], leverages multiple conditioning modalities—text,
570 audio, and video—within a single model. Another approach
571 in TV2A frameworks integrates ControlNet [54] to embed
572 video characteristics into text-to-audio synthesis, as demon-
573 strated by FoleyCrafter [55] and ReWAS [20]. FoleyCrafter
574 extracts frame-based clips as global features in IP-Adapter
575 [53] and trains a timestamp detector to identify sound effect
576 occurrences, integrating this information into ControlNet.

577 Our work is built on StableAudio Open [9] and
578 TangoFlux [16], designed for high-quality text-to-audio gen-
579 eration. StableAudio Open is a latent diffusion model that
580 generates stereo audio up to 47 seconds from text input, uti-
581 lizing a T5 text encoder for text processing and enabling
582 control over output length. In contrast, TangoFlux is based
583 on rectified flow, producing stereo audio up to 30 seconds
584 at 44.1 kHz, while leveraging a pretrained autoencoder from
585 StableAudio Open to enhance efficiency. Additionally, Tan-
586 goFlux incorporates CLAP-Ranked Preference Optimiza-
587 tion (CRPO) to generate and refine audio preference data.

588 8. Conclusion

589 In this work, we presented CAFA, a controllable Auto-
590 matic Foley Artist designed for the video-and-text-to-audio
591 task. Our model ensures high-quality audio synthesis while
592 maintaining both temporal and semantic alignment with the
593 input video. Guided by text prompts, it allows users to
594 incorporate details beyond what is present in the video or
595 even introduce new creative elements. This capability en-
596 hances flexibility in sound design beyond video-only Foley
597 models. By leveraging the modality adapter, our approach
598 achieves strong performance on a low computational bud-
599 get. Both objective metrics and human evaluations con-
600 firm its effectiveness in generating high-quality, contextu-
601 ally relevant audio. We believe that further advancements
602 in video feature extraction and T2A model refinement will
603 help address outstanding challenges, such as synthesizing
multiple audio sources simultaneously and capturing finer
motion details in video.

604

References

- [1] Vanessa Theme Ament. *The Foley grail: The art of performing sound for film, games, and animation*. Routledge, 2014.
- [2] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset, 2020.
- [3] Hong-You Chen, Zhengfeng Lai, Haotian Zhang, Xinze Wang, Marcin Eichner, Keen You, Meng Cao, Bowen Zhang, Yinfei Yang, and Zhe Gan. Contrastive localized language-image pre-training. *arXiv preprint arXiv:2410.02746*, 2024.
- [4] Peihao Chen, Yang Zhang, Mingkui Tan, Hongdong Xiao, Deng Huang, and Chuang Gan. Generating visually aligned sound from videos. *IEEE Transactions on Image Processing*, 29:8292–8302, 2020.
- [5] Ziyang Chen, Prem Seetharaman, Bryan Russell, Oriol Nieto, David Bourgin, Andrew Owens, and Justin Salamon. Video-guided foley sound generation with multimodal controls. *arXiv preprint arXiv:2411.17698*, 2024.
- [6] Ziyang Chen, Prem Seetharaman, Bryan Russell, Oriol Nieto, David Bourgin, Andrew Owens, and Justin Salamon. Video-guided foley sound generation with multimodal controls. 2025.
- [7] Ho Kei Cheng, Masato Ishii, Akio Hayakawa, Takashi Shibuya, Alexander Schwing, and Yuki Mitsufuji. Taming multimodal joint training for high-quality video-to-audio synthesis. *arXiv preprint arXiv:2412.15322*, 2024.
- [8] Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. Fast timing-conditioned latent audio diffusion. In *Forty-first International Conference on Machine Learning*, 2024.
- [9] Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Stable audio open. *arXiv preprint arXiv:2407.14358*, 2024.
- [10] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction guided latent diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 3590–3598, 2023.
- [11] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind one embedding space to bind them all. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15180–15190, 2023.
- [12] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [14] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981.
- [15] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. In *International Conference on Machine Learning*, pages 13916–13932. PMLR, 2023.
- [16] Chia-Yu Hung, Navonil Majumder, Zhifeng Kong, Ambuj Mehrish, Rafael Valle, Bryan Catanzaro, and Soujanya Poria. Tangoflux: Super fast and faithful text to audio generation with flow matching and clap-ranked preference optimization. *arXiv preprint arXiv:2412.21037*, 2024.
- [17] Vladimir Iashin and Esa Rahtu. Taming visually guided sound generation. *arXiv preprint arXiv:2110.08791*, 2021.
- [18] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Sparse in space and time: Audio-visual synchronisation with trainable selectors. *arXiv preprint arXiv:2210.07055*, 2022.
- [19] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Syncformer: Efficient synchronization from sparse cues. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5325–5329. IEEE, 2024.
- [20] Yujin Jeong, Yunji Kim, Sanghyuk Chun, and Jiyoung Lee. Read, watch and scream! sound generation from text and video. *arXiv preprint arXiv:2407.05551*, 2024.
- [21] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance: A metric for evaluating music enhancement algorithms, 2019.
- [22] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition, 2020.
- [23] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*, 2022.
- [24] Solomon Kullback and R. A. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951.
- [25] FirstName LastName. Supplemental material - cafa: a controllable automatic foley artist, 2025. Supplied as supplemental material supplemental_material.pdf.
- [26] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023.
- [27] Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation with self-supervised pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [28] Xiulong Liu, Kun Su, and Eli Shlizerman. Tell what you hear from what you see—video to audio generation through text. *arXiv preprint arXiv:2411.05679*, 2024.
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [30] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.
- [31] Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. *Advances in Neural Information Processing Systems*, 36:48855–48876, 2023.

- 717 [32] Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal,
718 Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria.
719 Tango 2: Aligning diffusion-based text-to-audio generations
720 through direct preference optimization. In *Proceedings of the*
721 *32nd ACM International Conference on Multimedia*, pages
722 564–572, 2024.
- 723 [33] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Repre-
724 sentation learning with contrastive predictive coding. *arXiv*
725 preprint arXiv:1807.03748, 2018.
- 726 [34] Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra,
727 Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi,
728 and Joao F Henriques. Keeping your eye on the ball: Tra-
729 jectory attention in video transformers. *Advances in neural*
730 *information processing systems*, 34:12493–12506, 2021.
- 731 [35] William Peebles and Saining Xie. Scalable diffusion models
732 with transformers. In *Proceedings of the IEEE/CVF interna-*
733 *tional conference on computer vision*, pages 4195–4205,
734 2023.
- 735 [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
736 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
737 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
738 transferable visual models from natural language supervi-
739 sion. In *International conference on machine learning*, pages
740 8748–8763. PMLR, 2021.
- 741 [37] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-
742 pher D Manning, Stefano Ermon, and Chelsea Finn. Direct
743 preference optimization: Your language model is secretly a
744 reward model. *Advances in Neural Information Processing*
745 *Systems*, 36:53728–53741, 2023.
- 746 [38] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee,
747 Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and
748 Peter J Liu. Exploring the limits of transfer learning with a
749 unified text-to-text transformer. *Journal of machine learning*
750 *research*, 21(140):1–67, 2020.
- 751 [39] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
752 Patrick Esser, and Björn Ommer. High-resolution image
753 synthesis with latent diffusion models. In *Proceedings of*
754 *the IEEE/CVF conference on computer vision and pattern*
755 *recognition*, pages 10684–10695, 2022.
- 756 [40] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-
757 net: Convolutional networks for biomedical image segmen-
758 tation. In *Medical image computing and computer-assisted*
759 *intervention-MICCAI 2015: 18th international conference,*
760 *Munich, Germany, October 5-9, 2015, proceedings, part III*
761 *18*, pages 234–241. Springer, 2015.
- 762 [41] Tim Salimans and Jonathan Ho. Progressive distillation
763 for fast sampling of diffusion models. *arXiv* preprint
764 arXiv:2202.00512, 2022.
- 765 [42] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki
766 Cheung, Alec Radford, and Xi Chen. Improved techniques
767 for training gans. *ArXiv*, abs/1606.03498, 2016.
- 768 [43] Roy Sheffer and Yossi Adi. I hear your true colors: Im-
769 age guided audio generation. In *ICASSP 2023-2023 IEEE*
770 *International Conference on Acoustics, Speech and Signal*
771 *Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- 772 [44] Jiaming Song, Chenlin Meng, and Stefano Ermon.
773 Denoising diffusion implicit models. *arXiv* preprint
774 arXiv:2010.02502, 2020.
- 775 [45] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Ab-
776 hishek Kumar, Stefano Ermon, and Ben Poole. Score-based
777 generative modeling through stochastic differential equa-
778 tions. *arXiv preprint arXiv:2011.13456*, 2020.
- 779 [46] Ilpo Viertola, Vladimir Iashin, and Esa Rahtu. Temporally
780 aligned audio for video with autoregression. *arXiv preprint*
781 arXiv:2409.13689, 2024.
- 782 [47] Heng Wang, Jianbo Ma, Santiago Pascual, Richard
783 Cartwright, and Weidong Cai. V2a-mapper: A lightweight
784 solution for vision-to-audio generation by connecting foun-
785 dation models. In *Proceedings of the AAAI Conference on*
786 *Artificial Intelligence*, pages 15492–15501, 2024.
- 787 [48] Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jiawei
788 Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao.
789 Frieren: Efficient video-to-audio generation network with
790 rectified flow matching. *Advances in Neural Information*
791 *Processing Systems*, 37:128118–128138, 2025.
- 792 [49] Yusong Wu*, Ke Chen*, Tianyu Zhang*, Yuchen Hui*, Tay-
793 lor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale con-
794 trastive language-audio pretraining with feature fusion and
795 keyword-to-caption augmentation. In *IEEE International*
796 *Conference on Acoustics, Speech and Signal Processing,*
797 *ICASSP*, 2023.
- 798 [50] Yazhou Xing, Yingqing He, Zeyue Tian, Xintao Wang, and
799 Qifeng Chen. Seeing and hearing: Open-domain visual-
800 audio generation with diffusion latent aligners. In *Proced-
801 ings of the IEEE/CVF Conference on Computer Vision and*
802 *Pattern Recognition*, pages 7151–7161, 2024.
- 803 [51] Manjie Xu, Chenxing Li, Xinyi Tu, Yong Ren, Rilin Chen,
804 Yu Gu, Wei Liang, and Dong Yu. Video-to-audio generation
805 with hidden alignment. *arXiv* preprint arXiv:2407.07464,
806 2024.
- 807 [52] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao
808 Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete
809 diffusion model for text-to-sound generation. *IEEE/ACM*
810 *Transactions on Audio, Speech, and Language Processing*,
811 31:1720–1733, 2023.
- 812 [53] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-
813 adapter: Text compatible image prompt adapter for text-to-
814 image diffusion models. *arXiv* preprint arXiv:2308.06721,
815 2023.
- 816 [54] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
817 conditional control to text-to-image diffusion models. In
818 *Proceedings of the IEEE/CVF international conference on*
819 *computer vision*, pages 3836–3847, 2023.
- 820 [55] Yiming Zhang, Yicheng Gu, Yanhong Zeng, Zhening Xing,
821 Yuancheng Wang, Zhizheng Wu, and Kai Chen. Foley-
822 crafter: Bring silent videos to life with lifelike and synchro-
823 nized sounds. *arXiv* preprint arXiv:2407.01494, 2024.